

In the claims:

1. (Currently Amended) In a content delivery system having m servers, $S' = \{S_1, \dots, S_m\}$, n active customers, $C' = \{C_1, \dots, C_n\}$, and g geographic locations, $G' = \{G_1, \dots, G_g\}$, wherein $sdel_k$ is a server delay of server S_k , $ndel_{j,k}$ is a network delay observed by customers in geographic location G_j while retrieving content from server S_k , p_j is a priority value for customer C_i , c_i is a total load of customer C_i , $u_{i,j}$ is a fraction of requests coming to customer C_i from region G_j , $a_{i,j,k}$ is a mapping representing a fraction of requests coming to customer C_i from region G_j that have been redirected to server S_k , and s_k represents a load capacity of server S_k , a method for distributing server loads, the method comprising the steps of:

representing an average prioritized observed response time as

$$AORT = \frac{\sum_{i=1}^n \sum_{j=1}^g \sum_{k=1}^m a_{i,j,k} \times u_{i,j} \times c_i \times p_i \times (sdel_k + ndel_{j,k})}{\sum_{i=1}^n c_i \times p_i}; \text{ and}$$

generating a mapping that assigns requests from customers to a particular server while minimizing $AORT$; and
assigning requests to the particular server based on the generated mapping.

2. (Currently Amended) The A-method as recited in claim 1, wherein the assigning further including the step of assigning all requests from all customers in all regions to a particular server such that, for each $C_i \in C'$, $G_g \in G'$, $\sum_{k=1}^m a_{i,j,k} = 1.0$.

3. (Currently Amended) The A-method as recited in claim 1, wherein the assigning step assigns further including the step of assigning requests to a particular server while ensuring that the load capacity of each server is not exceeded such that, for each for

$$\text{each } S_k \in S', \sum_{i=1}^n \sum_{j=1}^g a_{i,j,k} \times u_{i,j} \times c_i \leq s_k.$$

35 4. (Currently Amended) The A-method as recited in claim 1, wherein the
 36 assigning step includes further including the step of assigning requests to a particular server
 37 while balancing the load of each server to within a maximum allowed deviation from a
 38 balanced state Θ such that, for all pairs of servers, S_k and S_l ,
 39

$$40 \frac{\sum_{i=1}^n \sum_{j=1}^g a_{i,j,k} \times u_{i,j} \times c_i}{\sum_{i=1}^n \sum_{j=1}^g a_{i,j,l} \times u_{i,j} \times c_i} \leq (1 + \Theta) \times \frac{s_k}{s_l}.$$

41 5. (Currently Amended) The A-method as recited in claim 4, wherein if
 42 the content delivery system should add one or more servers, or remove one or more
 43 customers, the load of each server will not be redistributed unless the maximum allowed
 44 deviation from a balanced state Θ is exceeded.
 45

46 6. (Currently Amended) The A-method as recited in claim 1, further
 47 including the step of adding one or more customers to the content delivery system only if
 48 $AORT_{new} \leq (1 + \Phi) \times AORT_{old}$:
 49

50 wherein $AORT_{old}$ and $AORT_{new}$ are old and new values of $AORT$

$$51 \text{ defined as } AORT_{old} = \frac{\sum_{i=1}^n \sum_{j=1}^g \sum_{k=1}^m a_{i,j,k} \times u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})}{\sum_{i=1}^n c_i}$$

$$52 \text{ and } AORT_{new} = \frac{\sum_{i=1}^n \sum_{j=1}^g \sum_{k=1}^m a'_{i,j,k} \times u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})}{\sum_{i=1}^n c_i};$$

53 wherein $a'_{i,j,k}$ is a new mapping resulting from the addition of one or
 54 more customers; and
 55 wherein Φ is an allowable change in $AORT$ for existing clients.

1 7. (Currently Amended) The A-method as recited in claim 1, further
 2 including the step of using a linear constraint solver to generate the mapping.

3 8. (Currently Amended) The A-method as recited in claim 1, further
 4 including the step of using a non-linear constraint solver to generate the mapping.
 5 (Currently Amended) The A-method as recited in claim 1, further including the step of using
 6 a heuristic algorithm to generate the mapping, the heuristic algorithm comprising the step of
 7 assigning large $a_{i,j,k}$ values to small $u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})$ values to produce a smaller overall
 8 AORT value.

9 10. (Currently Amended) The A-method as recited in claim 9, the heuristic
 10 algorithm comprising the steps of:
 11 generating a plurality of sorted lists by sorting C_i values in increasing
 12 order of c_i , sorting $\langle C_i, G_j \rangle$ pairs in increasing order of $u_{i,j}$, sorting S_k values in increasing order
 13 of $sdel_k$, and sorting $\langle G_j, S_k \rangle$ pairs in increasing order of $ndel_{j,k}$;
 14 starting with a top-most, smallest value item in each list, identifying
 15 comparable smallest-value items from the other lists to generate a plurality of $\langle C_i, G_j, S_k \rangle$
 16 triples equivalent to the number of sorted lists;
 17 selecting from the plurality of $\langle C_i, G_j, S_k \rangle$ triples, the $\langle C_i, G_j, S_k \rangle$ triple
 18 with the smallest $u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})$ value;
 19 assigning to a server S_k of the selected $\langle C_i, G_j, S_k \rangle$ triple a remaining load
 20 from the $\langle C_i, G_j \rangle$ pair; and
 21 repeating the heuristic algorithm starting with generating the plurality
 22 of sorted lists, taking into account the changes in the values of the C_i values and the $\langle C_i, G_j \rangle$
 23 pairs as a result of the previous server assignment during each iteration, until the load from all
 24 $\langle C_i, G_j \rangle$ pairs has been assigned to a server S_k ;
 25 wherein if, during any iteration of the heuristic algorithm, the load
 26 capacity of the server S_k is not sufficient to handle the remaining load, the remaining load
 27 capacity of the server S_k is assigned to some of the load of the $\langle C_i, G_j \rangle$ pair, and an unassigned
 28 portion of the load from the $\langle C_i, G_j \rangle$ pair is reinserted into the iterative process.

015.683623.1

11. (Currently Amended) The A-method as recited in claim 10, the
 heuristic algorithm further including the steps of:
 generating a load-capacity prioritized sorted list by sorting C_i values in
 decreasing order of remaining load capacity s_k ;
 starting with a top-most, largest value item in the load-capacity
 prioritized list, identifying comparable smallest-value items from the other lists to generate a
 load-capacity prioritized $\langle C_i, G_j, S_k \rangle$ triple;
 considering the load-capacity prioritized $\langle C_i, G_j, S_k \rangle$ triple in the
 selection of the top-most $\langle C_i, G_j, S_k \rangle$ triple with the smallest $u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})$ value; and
 repeating the heuristic algorithm starting with generating the plurality
 of sorted lists, taking into account the changes in the values of the C_i values, the $\langle C_i, G_j \rangle$ pairs,
 and remaining load capacity as a result of the previous server assignment during each
 iteration, until the load from all $\langle C_i, G_j \rangle$ pairs has been assigned to a server S_k .

12. (Currently Amended) The A-method as recited in claim 10, the
 heuristic algorithm further including the steps of:
 generating a list of content-available $\langle C_i, S_k \rangle$ pairs in which the content
 of customer C_i is stored in server S_k ; and
 selecting the $\langle C_i, G_j, S_k \rangle$ triple with the smallest $u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})$
 value that is also part of the list of content-available $\langle C_i, S_k \rangle$ pairs;
 wherein if, during any iteration of the heuristic algorithm, there is no
 $\langle C_i, G_j, S_k \rangle$ triple that is also part of the list of content-available $\langle C_i, S_k \rangle$ pairs, a suitable
 $\langle C_i, G_j, S_k \rangle$ triple with the smallest $u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})$ value is chosen and the data of
 customer C_i is migrated to server S_k .

13. (Currently Amended) The A-method as recited in claim 12, the
 heuristic algorithm further including the step of generating a list of content-unavailable
 $\langle C_i, S_k \rangle$ pairs in increasing order of migration time penalty for which the content of customer
 C_i is not stored in server S_k ;

wherein if, during any iteration of the heuristic algorithm, there is no
 $\langle C_i, G_j, S_k \rangle$ triple that is also part of the list of content-available $\langle C_i, S_k \rangle$ pairs, a suitable
 $\langle C_i, G_j, S_k \rangle$ triple with the smallest combined $u_{ij} \times c_i \times (sdel_k + ndel_{j,k})$ value and $\langle C_i, S_k \rangle$ migration
time penalty is chosen and the data of customer C_i is migrated to server S_k .

14. (Currently Amended) The A-method as recited in claim 1, further
including the step of estimating $ndel_{j,k}$ for non-persistent connections using HTTP logs, the
step of estimating $ndel_{j,k}$ for non-persistent connections using HTTP logs comprising the steps
of:

computing an estimated round trip delay $\Delta r_{server, client}$ as $t_{con_req_rec}(i+1) -$
 $t_{resp_send_end}(i)$ from information stored in the HTTP logs, where $t_{con_req_rec}(i+1)$ represents a
time at which a connection request message is received by the server for an $(i+1)^{th}$ object, and
 $t_{resp_send_end}(i)$ represents a time at which the server stops sending an i^{th} object; and
computing the response time as

$(t_{resp_send_end} - t_{con_req_rec}) + 2 \times \frac{\Delta r_{server, client}}{2}$, where $t_{con_req_rec}$ represents a time at which a
connection request message is received by the server for an object, and $t_{resp_send_end}$ represents
a time at which the server stops sending the object.

15. (Currently Amended) The A-method as recited in claim 1, further
including the step of estimating $ndel_{j,k}$ for persistent connections using HTTP logs, the step of
estimating $ndel_{j,k}$ for persistent connections using HTTP logs comprising the steps of:

computing an estimated round trip delay $\Delta r_{server, client}$ as $t_{con_close_rec} -$
 $t_{resp_send_end}(last)$, where $t_{con_close_rec}$ represents a time at which the server receives a request to
close the persistent connection, and $t_{resp_send_end}(last)$ represents a time at which the server
stops sending a response for a last request; and
computing the response time as

$(t_{con_close_rec} - t_{con_req_rec}) + \frac{\Delta r_{server, client}}{2} - \frac{\Delta r_{server, client}}{2}$, where $t_{con_req_rec}$ represents a connection
request time.

1 16. (Withdrawn) A method for estimating a per-byte network delay α
 2 observed by a requesting entity in a geographic location while retrieving data from a server
 3 using TCP logs, the method comprising the steps of:
 4 computing an *immediate* round-trip delay Δr_{imm} between the server and
 5 the requesting entity as $t_{con_est} - t_{ack_send}$ from information stored in the TCP logs, where t_{con_est}
 6 is a time at which a connection with the server is established, and t_{ack_send} is a time at which
 7 the server sends an acknowledgement to a connection request message received from the
 8 requesting entity; and
 9 determining the per-byte network delay α using the computed
 10 immediate round-trip delay Δr_{imm} , an equation $t_{ack_receive}(k) - t_{data_send}(l) = \Delta r_{imm} + \alpha \times (k-l)$, and
 11 information stored in the TCP logs, where $t_{data_send}(l)$ is a time at which an l^{th} byte of the data
 12 has been sent by the server, and $t_{ack_receive}(k)$ is a time at which acknowledgement for a k^{th} ($k \geq l$)
 13 byte of the data is received by the server.

14 17 (Currently Amended) A method for estimating a response time
 15 observed by a requesting entity in a geographic location while retrieving objects from a server
 16 through a non-persistent connection using HTTP logs, the method comprising the steps of:
 17 computing an estimated round trip delay $\Delta r_{server,client}$ as $t_{con_req_rec}(i+1) -$
 18 $t_{resp_send_end}(i)$ from information stored in the HTTP logs, where $t_{con_req_rec}(i+1)$ represents a
 19 time at which a connection request message is received by the server from the requesting
 20 entity for an $(i+1)^{th}$ object, and $t_{resp_send_end}(i)$ represents a time at which the server stops
 21 sending an i^{th} object to the requesting entity; and
 22 computing the response time as
 23 $(t_{resp_send_end} - t_{con_req_rec}) + 2 \times \frac{\Delta r_{server,client}}{2}$, where $t_{con_req_rec}$ represents a time at which a
 24 connection request message is received by the server from the requesting entity for an object,
 25 and $t_{resp_send_end}$ represents a time at which the server stops sending the object to the requesting
 26 entity.

27 18. (Currently Amended) A method for estimating a response time
 28 observed by a requesting entity in a geographic location while retrieving objects from a server
 29 through a persistent connection using HTTP logs, the method comprising the steps of:

30 computing an estimated round trip delay $\Delta r_{server,client}$ as $t_{con_close_rec} -$
 31 $t_{resp_send_end}(last)$, where $t_{con_close_rec}$ represents a time at which the server receives a request to
 32 close the persistent connection, and $t_{resp_send_end}(last)$ represents a time at which the server
 33 stops sending a response for a last request; and
 34 computing the response time as

35 $(t_{con_close_rec} - t_{con_req_rec}) + \frac{\Delta r_{server,client}}{2} - \frac{\Delta r_{server,client}}{2}$, where $t_{con_req_rec}$ represents a connection
 36 request time.

37 19. (Currently Amended) In a content delivery network having m servers,
 38 $S' = \{S_1, \dots, S_m\}$, n active customers, $C' = \{C_1, \dots, C_n\}$, and g geographic locations, $G' =$
 39 $\{G_1, \dots, G_g\}$, a content delivery system for distributing server loads, the content delivery
 40 system comprising:

41 memory for storing a server delay $sdel_k$ of server S_k , a network delay
 42 $ndel_{j,k}$ observed by customers in geographic location G_j while retrieving content from server
 43 S_k , a priority value p_i for customer C_i , a total load c_i of customer C_i , a fraction of requests $u_{i,j}$
 44 coming to customer C_i from region G_j , a mapping $a_{i,j,k}$ representing a fraction of requests
 45 coming to customer C_i from region G_j that have been redirected to server S_k , and a load
 46 capacity s_k of server S_k ; and

47 a processor programmed for
 48 representing an average prioritized observed response time as

$$49 \quad AORT = \frac{\sum_{i=1}^n \sum_{j=1}^g \sum_{k=1}^m a_{i,j,k} \times u_{i,j} \times c_i \times p_i \times (sdel_k + ndel_{j,k})}{\sum_{i=1}^n c_i \times p_i}, \text{ and}$$

50 generating a mapping that assigns requests from customers to a
 51 particular server while minimizing $AORT$,

52 the processor further programmed for assigning requests to a
 53 particular server based on said mapping.

Atty. Dkt. No. 073303-0116

20. (Currently Amended) The A-system as recited in claim 19, the processor further programmed for assigning all requests from all customers in all regions to a particular server such that, for each $C_i \in C'$, $G_g \in G'$, $\sum_{k=1}^m a_{i,j,k} = 1.0$.

21. (Currently Amended) The A-system as recited in claim 19, the processor further programmed for assigning requests to a particular server while ensuring that the load capacity of each server is not exceeded such that, for each for each $S_k \in S'$,

$$\sum_{i=1}^n \sum_{j=1}^g a_{i,j,k} \times u_{i,j} \times c_i \leq s_k.$$

22. (Currently Amended) The A-system as recited in claim 19, the processor further programmed for assigning requests to a particular server while balancing the load of each server to within a maximum allowed deviation from a balanced state Θ such that,

$$\text{for all pairs of servers, } S_k \text{ and } S_l, \frac{\sum_{i=1}^n \sum_{j=1}^g a_{i,j,k} \times u_{i,j} \times c_i}{\sum_{i=1}^n \sum_{j=1}^g a_{i,j,l} \times u_{i,j} \times c_i} \leq (1 + \Theta) \times \frac{s_k}{s_l}.$$

23. (Currently Amended) The A-system as recited in claim 22, wherein if the content delivery system should add one or more servers, or remove one or more customers, the processor is further programmed for not redistributing the load of each server unless the maximum allowed deviation from a balanced state Θ is exceeded.

24. (Currently Amended) The A-system as recited in claim 19, wherein the processor is programmed for allowing one or more customers to be added to the content delivery system only if $AORT_{new} \leq (1 + \Phi) \times AORT_{old}$:

Atty. Dkt. No. 073303-0116

wherein $AORT_{old}$ and $AORT_{new}$ are old and new values of $AORT$

$$\text{defined as } AORT_{old} = \frac{\sum_{i=1}^n \sum_{j=1}^g \sum_{k=1}^m a_{i,j,k} \times u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})}{\sum_{i=1}^n c_i}$$

$$\text{and } AORT_{new} = \frac{\sum_{i=1}^n \sum_{j=1}^g \sum_{k=1}^m a'_{i,j,k} \times u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})}{\sum_{i=1}^n c_i};$$

wherein $a'_{i,j,k}$ is a new mapping resulting from the addition of one or more customers; and

wherein Φ is an allowable change in $AORT$ for existing clients.

25. (Currently Amended) The A-system as recited in claim 19, the processor further programmed for using a linear constraint solver to generate the mapping.

26. (Currently Amended) The A-system as recited in claim 19, the processor further programmed for using a non-linear constraint solver to generate the mapping.

27. (Currently Amended) The A-system as recited in claim 19, the processor further programmed for generating the mapping by assigning large $a_{i,j,k}$ values to small $u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})$ values to produce a smaller overall $AORT$ value.

28. (Currently Amended) The A-system as recited in claim 27, the processor further programmed for:

Atty. Dkt. No. 073303-0116

generating a plurality of sorted lists by sorting C_i values in increasing order of c_i , sorting $\langle C_i, G_j \rangle$ pairs in increasing order of $u_{i,j}$, sorting S_k values in increasing order of $sdel_k$, and sorting $\langle G_j, S_k \rangle$ pairs in increasing order of $ndel_{j,k}$;

starting with a top-most, smallest value item in each list, identifying comparable smallest-value items from the other lists to generate a plurality of $\langle C_i, G_j, S_k \rangle$ triples equivalent to the number of sorted lists;

selecting from the plurality of $\langle C_i, G_j, S_k \rangle$ triples, the $\langle C_i, G_j, S_k \rangle$ triple with the smallest $u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})$ value;

assigning to a server S_k of the selected $\langle C_i, G_j, S_k \rangle$ triple a remaining load from the $\langle C_i, G_j \rangle$ pair; and

repeating the heuristic algorithm starting with generating the plurality of sorted lists, taking into account the changes in the values of the C_i values and the $\langle C_i, G_j \rangle$ pairs as a result of the previous server assignment during each iteration, until the load from all $\langle C_i, G_j \rangle$ pairs has been assigned to a server S_k ;

wherein if, during any iteration of the heuristic algorithm, the load capacity of the server S_k is not sufficient to handle the remaining load, the remaining load capacity of the server S_k is assigned to some of the load of the $\langle C_i, G_j \rangle$ pair, and an unassigned portion of the load from the $\langle C_i, G_j \rangle$ pair is reinserted into the iterative process.

29. (Currently Amended) The A-system as recited in claim 28, the processor further programmed for:

generating a load-capacity prioritized sorted list by sorting C_i values in decreasing order of remaining load capacity s_k ;

starting with a top-most, largest value item in the load-capacity prioritized list, identifying comparable smallest-value items from the other lists to generate a load-capacity prioritized $\langle C_i, G_j, S_k \rangle$ triple;

considering the load-capacity prioritized $\langle C_i, G_j, S_k \rangle$ triple in the selection of the top-most $\langle C_i, G_j, S_k \rangle$ triple with the smallest $u_{i,j} \times c_i \times (sdel_k + ndel_{j,k})$ value; and

repeating the heuristic algorithm starting with generating the plurality of sorted lists, taking into account the changes in the values of the C_i values, the $\langle C_i, G_j \rangle$ pairs,

Atty. Dkt. No. 073303-0116

and remaining load capacity as a result of the previous server assignment during each iteration, until the load from all $\langle C_i, G_j \rangle$ pairs has been assigned to a server S_k .

30. (Currently Amended) The A-system as recited in claim 28, the processor further programmed for:

- generating a list of content-available $\langle C_i, S_k \rangle$ pairs in which the content of customer C_i is stored in server S_k ; and
- selecting the $\langle C_i, G_j, S_k \rangle$ triple with the smallest $u_{ij} \times c_i \times (sdel_k + ndel_{j,k})$ value that is also part of the list of content-available $\langle C_i, S_k \rangle$ pairs;
- wherein if, during any iteration of the heuristic algorithm, there is no $\langle C_i, G_j, S_k \rangle$ triple that is also part of the list of content-available $\langle C_i, S_k \rangle$ pairs, a suitable $\langle C_i, G_j, S_k \rangle$ triple with the smallest $u_{ij} \times c_i \times (sdel_k + ndel_{j,k})$ value is chosen and the data of customer C_i is migrated to server S_k .

31. (Currently Amended) The A-system as recited in claim 30, the processor further programmed for generating a list of content-unavailable $\langle C_i, S_k \rangle$ pairs in increasing order of migration time penalty for which the content of customer C_i is not stored in server S_k ;

- wherein if, during any iteration of the heuristic algorithm, there is no $\langle C_i, G_j, S_k \rangle$ triple that is also part of the list of content-available $\langle C_i, S_k \rangle$ pairs, the processor is further programmed for selecting a suitable $\langle C_i, G_j, S_k \rangle$ triple with the smallest combined $u_{ij} \times c_i \times (sdel_k + ndel_{j,k})$ value and $\langle C_i, S_k \rangle$ migration time penalty, and migrating the data of customer C_i to server S_k .

32. (Currently Amended) The A-system as recited in claim 19, the processor further programmed for estimating $ndel_{j,k}$ for non-persistent connections using HTTP logs by:

- computing an estimated round trip delay $\Delta t_{server, client}$ as $t_{con_req_rec}(i+1) - t_{resp_send_end}(i)$ from information stored in the HTTP logs, where $t_{con_req_rec}(i+1)$ represents a time at which a connection request message is received by the server for an $(i+1)^{th}$ object, and $t_{resp_send_end}(i)$ represents a time at which the server stops sending an i^{th} object; and

Atty. Dkt. No. 073303-0116

computing the response time as

$(t_{resp_send_end} - t_{con_req_rec}) + 2 \times \frac{\Delta r_{server,client}}{2}$, where $t_{con_req_rec}$ represents a time at which a connection request message is received by the server for an object, and $t_{resp_send_end}$ represents a time at which the server stops sending the object.

33. (Currently Amended) The A-system as recited in claim 19, the processor further programmed for estimating $ndel_{j,k}$ for persistent connections using HTTP logs by:

computing an estimated round trip delay $\Delta r_{server,client}$ as $t_{con_close_rec} - t_{resp_send_end}(last)$, where $t_{con_close_rec}$ represents a time at which the server receives a request to close the persistent connection, and $t_{resp_send_end}(last)$ represents a time at which the server stops sending a response for a last request; and

computing the response time as

$(t_{con_close_rec} - t_{con_req_rec}) + \frac{\Delta r_{server,client}}{2} - \frac{\Delta r_{server,client}}{2}$, where $t_{con_req_rec}$ represents a connection request time.